

Algorithm A.11 Conjugate gradient ascent

```

Procedure Conjugate-Gradient-Ascent (
   $\theta^1$ , // Initial starting point
   $f_{\text{obj}}$ , // Function to be optimized
   $\delta$  // Convergence threshold
)
1   $t \leftarrow 1$ 
2   $g^0 \leftarrow \mathbf{1}$ 
3   $h^0 \leftarrow \mathbf{0}$ 
4  do
5     $g^t \leftarrow \nabla f_{\text{obj}}(\theta^t)$ 
6     $\gamma^t \leftarrow \frac{(g^t - g^{t-1})^T g^t}{(g^{t-1})^T g^{t-1}}$ 
7     $h^t \leftarrow g^t + \gamma^t h^{t-1}$ 
8    Choose  $\eta^t$  by line search along the line  $\theta_t + \eta h^t$ 
9     $\theta^{t+1} \leftarrow \theta^t + \eta^t h^t$ 
10    $t \leftarrow t + 1$ 
11  while  $\|\theta^t - \theta^{t-1}\| > \delta$ 
12  return ( $\theta^t$ )

```

A.5.3 Constrained Optimization

In appendix A.5.1, we considered the problem of optimizing a continuous function over its entire domain (see also appendix A.5.2). In many cases, however, we have certain constraints that the desired solution must satisfy. Thus, we have to optimize the function within a constrained space. We now review some basic methods that address this problem of *constrained optimization*.

constrained
optimization

Example A.5

Suppose we want to find the maximum entropy distribution over a variable X , with $\text{Val}(X) = \{x^1, \dots, x^K\}$. Consider the entropy of X :

$$H(X) = - \sum_{k=1}^K P(x^k) \log P(x^k).$$

We can maximize this function using the gradient method by treating each $P(x^k)$ as a separate parameter θ_k . We compute the gradient of $H_P(X)$ with respect to each of these parameters:

$$\frac{\partial}{\partial \theta_k} H(X) = -\log(\theta_k) - 1.$$

Setting this partial derivative to 0, we get that $\log(\theta_k) = -1$, and thus $\theta_k = 1/2$. This solution seems fine until we realize that the numbers do not sum up to 1, and hence our solution does not define a probability distribution!

The flaw in our analysis is that we want to maximize the entropy subject to a constraint on the parameters, namely, $\sum_k \theta_k = 1$. In addition, we also remember that we need to require that $\theta_k \geq 0$. In this case we see that the gradient drives the solution away from from 0 ($-\log(\theta_k) \rightarrow \infty$ as $\theta_k \rightarrow 0$), and thus we do not need to enforce this constraint actively. ■